

Object Instance Retrieval in Assistive Robotics: Leveraging Fine-Tuned SimSiam with Multi-View Images Based on 3D Semantic Map

Taichi Sakguchi, Akira Taniguchi, Yoshinobu Hagiwara,
Lotfi El Hafi, Shoichi Hasegawa, Tadahiro Taniguchi

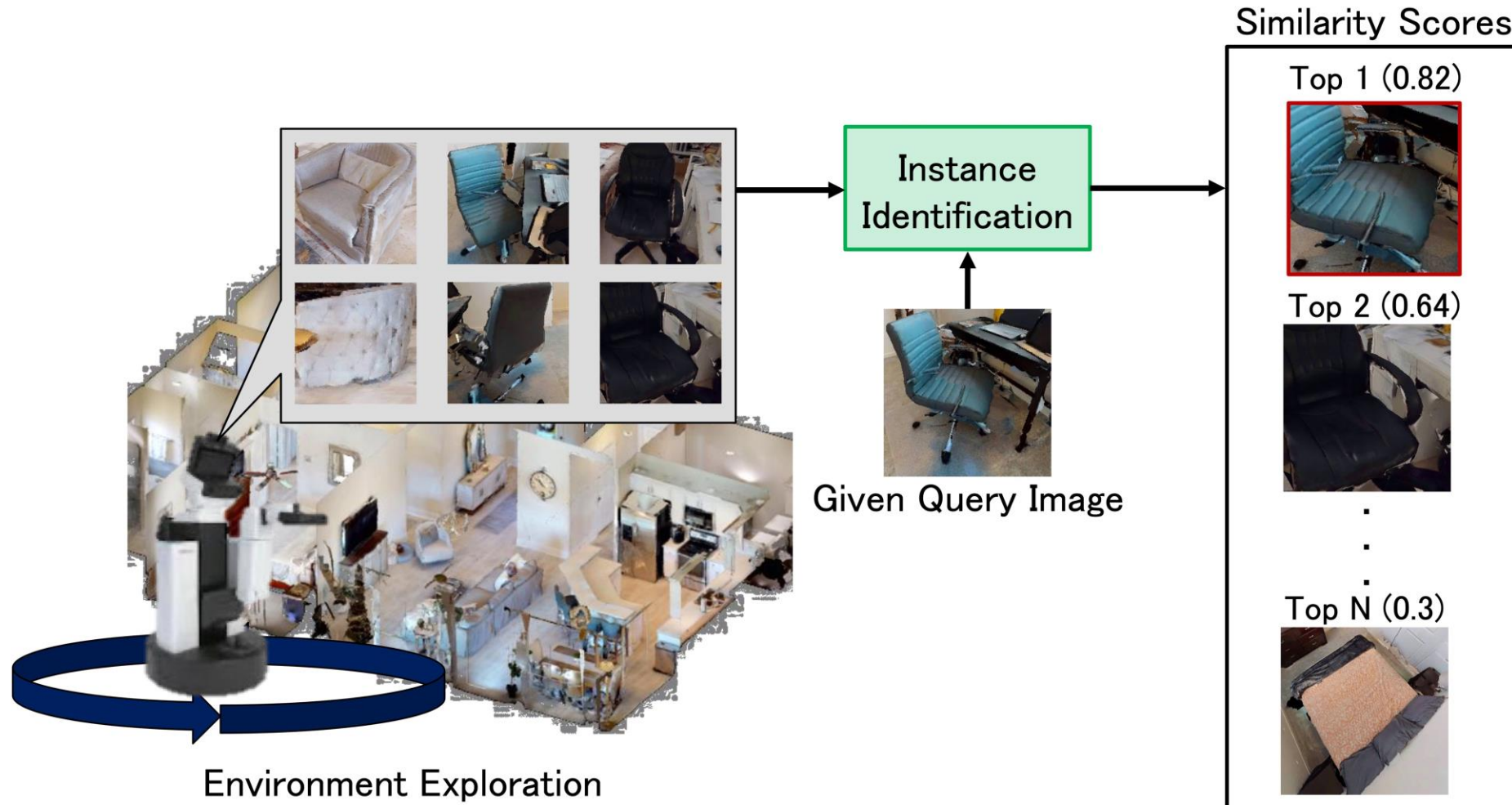
Ritsumeikan University, Japan



Task: Instance-specific Image Goal Navigation (InstanceImageNav)

InstanceImageNav:

The task of finding instances identical to a query image in an environment with multiple instances of the same object class, each with a different appearance.



Problem Statements:

Problem 1:

To solve InstanceImageNav, robots need to identify different instances of the same class of objects.

Contrastive language-image pre-training (CLIP), which was used in previous research on object search, is not suitable for fine-grained identification tasks.

Problem2 :

When a robot explores the environment and observes an object, it observes the object from various 3D viewpoints.



In this case, the appearance could be significantly different and could decrease the success rate of InstanceImageNav

Hypothesizes:

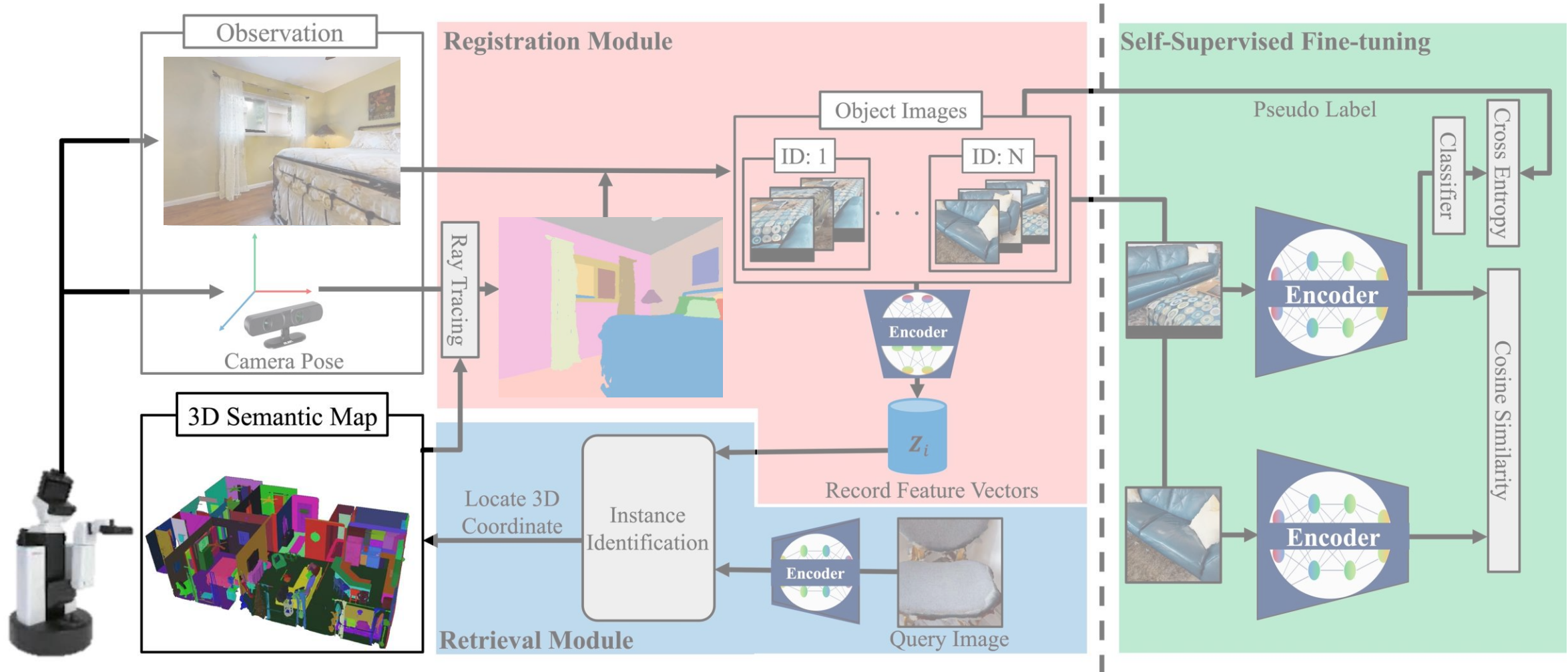
1. Models from contrastive learning between images could be superior to CLIP for instance-level object identification.
2. Increasing image similarity across various 3d viewpoints in self-supervised learning could improve instance retrieval accuracy.



Proposal: SimView

SimView:

Leveraging multi-view images based on a 3D semantic map of the environment and contrastive learning to train an instance identification model on-site



Simulation experiments: Quantitative result

Metrics: mean Average Precision

Comparison Method

- Unimodal Contrastive Learning Method
 - SimSiam, DINOv2, SimCLR
- Multimodal Contrastive Learning Method
 - CLIP

| Method | Arch. | Env. 1 | Env. 2 | Env. 3 | Env. 4 | Env. 5 | Env. 6 | Env. 7 | Env. 8 | Env. 9 | Avg. |
|---------|----------|--------------------|--------------------|--------------------|--------------------|--------------------|--------------------|--------------------|--------------------|--------------------|--------------------|
| SimView | ResNet50 | <u>0.79</u> | <u>0.67</u> | <u>0.68</u> | <u>0.91</u> | <u>0.68</u> | <u>0.45</u> | <u>0.72</u> | <u>0.74</u> | <u>0.78</u> | <u>0.72</u> |
| SimSiam | ResNet50 | <u>0.70</u> | <u>0.58</u> | <u>0.56</u> | 0.80 | 0.65 | 0.41 | 0.63 | <u>0.68</u> | <u>0.75</u> | <u>0.64</u> |
| DINOv2 | ViT-B/14 | 0.51 | 0.48 | 0.45 | 0.66 | <u>0.71</u> | 0.41 | 0.55 | 0.44 | 0.61 | 0.54 |
| SimCLR | ResNet50 | 0.65 | 0.56 | 0.52 | <u>0.83</u> | 0.64 | <u>0.45</u> | 0.66 | 0.62 | <u>0.75</u> | 0.63 |
| CLIP | ResNet50 | 0.51 | 0.36 | 0.35 | 0.51 | 0.48 | 0.44 | 0.42 | 0.44 | 0.59 | 0.46 |
| CLIP | ViT-B/16 | 0.56 | 0.38 | 0.37 | 0.47 | 0.48 | 0.37 | 0.38 | 0.35 | 0.48 | 0.43 |